

# A Machine Learning Framework for Predicting Drug-drug Interactions

Suyu Mei (✉ [061021053@fudan.edu.cn](mailto:061021053@fudan.edu.cn))

Shenyang Normal University

Kun Zhang

Xavier University of Louisiana

---

## Research Article

**Keywords:** drug-drug interaction, drug-target interaction, machine learning, protein-protein interaction networks, signaling pathways

**Posted Date:** May 11th, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-503867/v1>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

# 1 **A machine learning framework for predicting drug-drug** 2 **interactions**

3 Suyu Mei<sup>\*1</sup> and Kun Zhang<sup>\*2</sup>

4 <sup>1</sup>Software College, Shenyang Normal University, Shenyang 110034, China

5 <sup>2</sup>Bioinformatics Core of Xavier RCMC Center for Cancer Research, Department of  
6 Computer Science, Xavier University of Louisiana, New Orleans, LA 70125, US \*

7 Corresponding author: Suyu Mei (e-mail: [meisygle@gmail.com](mailto:meisygle@gmail.com)); Kun Zhang (e-mail:  
8 [kunzhang@xula.edu](mailto:kunzhang@xula.edu)).

## 9 **Abstract**

10 Understanding drug-drug interaction is an essential step to reduce the risk of adverse  
11 drug events before clinical drug co-prescription. Existing methods commonly integrate  
12 multiple heterogeneous data sources to increase model performance but result in a high  
13 model complexity. To elucidate the molecular mechanisms behind drug-drug  
14 interactions and reserve rational biological interpretability is a major concern in  
15 computational modeling. In this study, we propose a simple representation of drug  
16 target profiles to depict drug pairs, based on which an  $l_2$ -regularized logistic regression  
17 model is built to predict drug-drug interactions. In addition, we develop several  
18 statistical metrics to measure the communication intensity, interaction efficacy and  
19 action range between two drugs in the context of human protein-protein interaction  
20 networks and signaling pathways. Cross validation and independent test show that the  
21 simple feature representation via drug target profiles is effective to predict drug-drug  
22 interactions and outperforms the existing data integration methods. Statistical results  
23 show that two drugs easily interact when they target common genes, or their target  
24 genes communicate with each other via short paths in protein-protein interaction  
25 networks or through cross-talks between signaling pathways. The unravelled

26 mechanisms provide biological insights into potential pharmacological risks of known  
27 drug-drug interactions and drug target genes.

28 **Keywords** drug-drug interaction; drug-target interaction; machine learning; protein-  
29 protein interaction networks; signaling pathways

## 30 **Introduction**

31 Drug-drug interactions (DDIs) have been recognized as a major cause of adverse drug  
32 reactions (ADRs) that leads to increasing healthcare costs [1]. DDIs may occur  
33 synergistically or antagonistically and potentially result in adverse side effects and  
34 toxicities when a patient takes more than one drug concurrently [2]. In many cases,  
35 DDIs are hardly detected during clinical trial phase and arbitrary co-prescription of  
36 drugs without prior knowledge potentially poses a serious threat to patient health and  
37 life [3]. To reduce the risk of potential side effects, we need to check in advance whether  
38 two co-prescribed drugs interact. DDIs could be identified via *in vitro*, *in vivo*  
39 experiments and *in silico* computational methods. The two former approaches are very  
40 costly and are in some cases impossible to be carried out, because DDIs could elicit  
41 serious side effects in *in vitro* or *in vivo* experiments [4]. With the advances of  
42 pharmacogenomics, recent years have witnessed much progress in developing *in silico*  
43 computational methods to predict DDIs and their effects.

44 Existing computational methods could be roughly classified into three categories, i.e.  
45 similarity-based methods [5-8], networks-based methods [9-13] and machine learning  
46 methods [14-22]. Similarity-based methods directly calculate similarity scores between  
47 drug profiles to infer DDIs. Vilar et al. [5] review several drug profiles that have been  
48 used to infer drug repurposing, adverse effects detection and drug–drug interactions,  
49 e.g. pharmaceutical profiles, gene expression profiles and phenome profiles. Among

50 these profiles, drug structural profile could be well interpreted on the assumption that  
51 structurally similar drugs tend to target related genes to produce some biological effects  
52 [6]. Besides drug profiles, similarity metric is the second major concern in similarity-  
53 based methods. Ferdousi et al. [7] compare a dozen of similarity metrics (e.g. inner  
54 product, Jaccard similarity, Russell-Rao similarity, Tanimoto coefficient) on drug  
55 target profiles and select the optimum measure to infer DDIs. Similarity-based methods  
56 have the merit of simple and intuitive interpretation, but the major disadvantage is that  
57 these methods cannot resist noise because the thresholding of similarity scores would  
58 be seriously affected by potential false DDIs. Networks-based methods could be further  
59 divided into two categories, i.e. drug similarity networks-based methods [9-11] and  
60 protein-protein interaction (PPI) networks-based methods [12, 13]. The former category  
61 of methods first constructs drug-drug similarity networks via a variety of drug similarity  
62 metrics and then predicts novel links/DDIs via networks inference, e.g. matrix  
63 factorization [9, 10], block coordinate descent optimization [11]. These methods,  
64 together with the similarity-based methods [5-8], focus on using drug structural  
65 similarities to infer DDIs. Actually, drug-drug interactions often refer to the biological  
66 events that two co-prescribed drugs influence or alter each other's therapeutic effects  
67 when they take actions on common genes or associated signaling pathways [7]. In this  
68 sense, mechanism information is more rational to infer DDIs than drug structural  
69 similarity, especially for the prediction of interactions between two structurally  
70 dissimilar drugs. Comparatively, the PPI networks-based methods [12, 13] could detect  
71 any pair of drugs that simultaneously act on associated genes to produce unexpected  
72 effects and thus gain a large coverage of DDIs with less bias. Park et al. [12] assume  
73 that DDIs are caused by close interference on common genes or associated genes within  
74 the same pathways as well as distant interference through cross-talking pathways, based

75 on which to capture distant interference via random walk algorithm on PPI networks.  
76 Huang et al. [13] also consider drug actions in the context of PPI networks and define  
77 the target genes together with their neighboring genes in PPI networks as a target-  
78 centred system for each drug, based on which the authors propose a S-score to measure  
79 the similarity between two drugs' target-centered systems. PPI networks-based  
80 methods have the merit of capturing the mechanism of drug actions but are restricted  
81 by the incompleteness of physical PPI networks.

82 In the last decades, machine learning has attracted increasing attention in inferring  
83 drug-drug interactions [14-22]. Most of the machine learning methods focus on data  
84 integration to increase the accuracy of DDI prediction. Data integration refers to  
85 capturing multiple aspects of a single data source or combining more than one data  
86 source. For instance, Dhimi et al. [14] explore heterogeneous similarities between drug  
87 structures from SMILES representation of drugs (i.e. molecular feature similarity,  
88 string similarity, molecular fingerprint similarity and molecular access system), while  
89 the other methods focus on integrating multiple data sources. Among the features, drug  
90 chemical structural data are the most frequently used [14-21] and are often combined  
91 with other data to predict DDIs, e.g. drug adverse drug reactions (ADR) [15-17, 20-21],  
92 target similarity [15-17, 19-21], PPI networks [20, 21] and signaling pathways [16].  
93 From methodological point of view, the existing machine learning approaches used to  
94 implement data integration majorly include ensemble learning [15, 16], kernel method  
95 [14, 17] and deep learning [18, 19]. Data integration is effective to combine multiple  
96 data sources or to capture multiple aspects of a single data source. Nevertheless, the  
97 performance improvement is often achieved at the cost of model complexity. The  
98 existing studies do not reveal which feature information is dispensable and contributes  
99 most to DDI prediction and which features are actually less informative. More

100 importantly, data integration imposes intensively demanding constraints on data that  
101 are potentially not simultaneously available (e.g. drug and target structures).  
102 Comparatively, single data source could reduce model complexity and achieve good  
103 interpretability. For instance, Karim et al. [22] use deep neural networks to  
104 automatically learn feature representation from the available DDI networks alone to  
105 predict novel DDIs. In addition, some inherent mechanisms surely exist behind drug-  
106 drug interactions, e.g. common drug-targeted genes and cross-talks between drug-  
107 targeted signaling pathways [23]. The underlying DDI mechanisms are often ignored  
108 by the existing machine learning methods that focus on drug molecular information  
109 itself.

110 In this study, we propose a DDI mechanisms based computational framework  
111 without the need of drug structural or ADR similarity information. In this framework,  
112 we assume that a drug alters other drugs' therapeutic effects through the associations  
113 between their targeted genes or signaling pathways and drug structural similarity is not  
114 indispensable for drug-drug interactions. We depict each drug using a drug target  
115 profile derived from DrugBank [24]. The profile is actually a binary vector with each  
116 element indicating the presence or absence of a gene. The target profiles of two drugs  
117 are simply combined as feature vector to depict the drug pair. We extract the known  
118 DDIs from DrugBank [24] as positive training data and randomly sample drug pairs as  
119 negative training data to train a  $l_2$ -regularized logistic regression model. The model is  
120 evaluated via cross validation on the training data and independent test on external data  
121 from the comprehensive database [25]. Lastly, we further conduct mechanism analyses  
122 of the known and the predicted DDIs for further biomedical research.

## 123 **Data and methods**

### 124 **Data**

125 The known drug-drug interactions (DDI) and drug-target interactions (DTI) are  
126 extracted from DrugBank [24]. As we use drug target profile to depict drugs and  
127 represent drug pairs, only the drugs that target at least one human gene are studied in  
128 this work. As such, we extract 6066 drugs and 2940 human target genes from DrugBank  
129 [24]. As results, we totally obtain 915,413 DDIs as the positive training data and obtain  
130 23,169 drug-gene interaction pairs to construct feature vectors representing drug pairs.  
131 From the 6066 drugs, we randomly sample 915,413 drug pairs disjoint with the positive  
132 training data as the negative training data.

133 The comprehensive database [25] has curated a large number of DDIs from  
134 experiments and text mining. After removing the DDIs that already exist in DrugBank  
135 [24], we totally obtain 13 independent datasets, e.g. 8188, 2003 and 37 DDIs from  
136 KEGG [26], OSCAR [27] (<https://oscar-emr.com/>) and VA NDF-RT [28], respectively.  
137 These datasets are used as positive independent test data. To estimate the risk of model  
138 bias, we randomly sample 8188 drug pairs as negative independent test data, which are  
139 disjoint with the training data and the positive independent test data.

140 To study the underlying molecular mechanisms associated with drug-drug  
141 interactions, we construct human physical protein-protein interaction (PPI) networks  
142 from HPRD [29], BioGRID [30], IntAct [31] and HitPredict [32]. We totally obtain  
143 171,249 human physical PPIs. In addition, we obtain 27 human immune signaling  
144 pathways from NetPath [33], in which IL1~IL11 are merged into one single pathway  
145 for simplicity. From Reactome [34], 1846 human signaling pathways are extracted.

## 146 Drug target profile-based feature construction

147 In this study, we investigate drug-drug interactions from the aspect of genes that two  
148 drugs act on. For each drug  $d_i$  in the DDI-associated drug set  $D$ , the human gene set  
149 targeted by  $d_i$  is assumed to be  $G_{d_i}$ . The entire target gene set is defined as follows.

$$150 \quad G = \cup_{d_i \in D} G_{d_i} \quad (1)$$

151 For each drug  $d_i$ , its target profile is formally defined as follows.

$$152 \quad V_{d_i}[g] = \begin{cases} 1, & g \in G_{d_i} \wedge g \in G \\ 0, & g \notin G_{d_i} \wedge g \in G \end{cases} \quad (2)$$

153 Then the feature vector of a drug pair  $(d_i, d_j)$  is represented by combining the target  
154 profile of  $d_i$  and  $d_j$  as follows.

$$155 \quad V_{(d_i, d_j)}[g] = V_{d_i}[g] + V_{d_j}[g], g \in G \quad (3)$$

156 The genes  $g \notin G$  are discarded. The simple feature representation as described by  
157 Formula (3) could intuitively reveal the co-occurrence of common target genes together  
158 with unique presence of some target genes between two drugs.

## 159 L<sub>2</sub>-regularized logistic regression as base learner

160 We adopt l<sub>2</sub>-regularized logistic regression [35] as base learner because it could fast fit  
161 large training data and penalize potential overtraining. Given training data  $x$  and labels  
162  $y$  consisting of instance-label pairs  $(x_i, y_i), i = 1, 2, \dots, l; x_i \in R^n; y_i \in \{-1, +1\}$ , the  
163 decision function of logistic regression is defined as  $f(x) = \frac{1}{1 + \exp(-y\omega^T x)}$ . L<sub>2</sub>-  
164 regularized logistic regression derives weight vector  $\omega$  via solving the following  
165 prime optimization problem.

$$166 \quad \min_{\omega} \frac{1}{2} \omega^T \omega + C \sum_{i=1}^l \log(1 + e^{-y_i \omega^T x_i}) \quad (4)$$

167 where  $C$  denotes the penalty parameter/regularizer and the second term penalizes  
 168 potential noise/outlier or overtraining. The optimization problem (4) is solved via its  
 169 dual form.

$$170 \min_{\alpha} \frac{1}{2} \alpha^T Q \alpha + \sum_{i:\alpha_i > 0} \alpha_i \log \alpha_i + \sum_{i:\alpha_i < C} (C - \alpha_i) \log (C - \alpha_i) - \sum_i^l C \log C \quad (5)$$

$$171 \text{ s. t. } 0 \leq \alpha_i \leq C, i = 1, \dots, l$$

172 where  $\alpha_i$  denotes Lagrangian operator and  $Q_{ij} = y_i y_j x_i^T x_j$ . To simplify the  
 173 parameter tuning, the regularizer  $C$  as defined in Formula (4) is chosen within the set  
 174  $\{2^i \mid -16 \leq i \leq 16, i \in I\}$ , where  $I$  denotes the integer set.

## 175 **Model evaluation and DDI mechanism metrics**

### 176 *Binary classification metrics*

177 We adopt five frequently-used performance metrics to evaluate the model performance,  
 178 i.e. Receiver Operating Characteristic curve AUC (ROC-AUC), sensitivity (SE),  
 179 precision (PR), Matthews correlation coefficient (MCC), Accuracy and F1 score. ROC-  
 180 AUC is calculated from the decision values produced by  $f(x)$ . The other metrics are  
 181 calculated via a confusion matrix  $M$ , whose element  $M_{i,j}$  records the counts that class  
 182  $i$  are classified to class  $j$ . From  $M$ , several intermediate variables are defined first as  
 183 Formula (6). Then the label-specific  $PR_l$ ,  $SE_l$  and  $MCC_l$  for each label are further  
 184 defined by Formula (7). The overall accuracy and MCC are defined by Formula (8).

$$186 p_l = M_{l,l}, q_l = \sum_{i=1, i \neq l}^L \sum_{j=1, j \neq l}^L M_{i,j}, r_l = \sum_{i=1, i \neq l}^L M_{i,l}, s_l = \sum_{j=1, j \neq l}^L M_{l,j}$$

$$185 p = \sum_{l=1}^L p_l, q = \sum_{l=1}^L q_l, r = \sum_{l=1}^L r_l, s = \sum_{l=1}^L s_l \quad (6)$$

$$187 PR_l = \frac{p_l}{p_l + r_l}, l = 1, 2, \dots, L$$

$$188 SE_l = \frac{p_l}{p_l + s_l}, l = 1, 2, \dots, L$$

189  $MCC_l = \frac{(p_l q_l - r_l s_l)}{\sqrt{(p_l + r_l)(p_l + s_l)(q_l + r_l)(q_l + s_l)}}, l = 1, 2, \dots, L$  (7)

191  $Acc = \frac{\sum_{l=1}^L M_{l,l}}{\sum_{i=1}^L \sum_{j=1}^L M_{i,j}}$

190  $MCC = \frac{(pq - rs)}{\sqrt{(p+r)(p+s)(q+r)(q+s)}}$  (8)

192 where  $L$  denotes the number of labels and assumes 2 in this study. F1 score is defined  
193 as follows.

194  $F1 \text{ score} = \frac{2 \times PR_l \times SE_l}{PR_l + SE_l}, l = 1 \text{ denotes the positive class}$  (9)

195 **DDI mechanism metrics**

196 Targeting common genes is one important way that a drug alters other drugs'  
197 therapeutic effects to result in drug-drug interactions. We use Jaccard index to measure  
198 the intensity that two drugs interact. For a drug pair  $(d_i, d_j)$ , the Jaccard index between  
199 them is defined as follows.

200  $Jaccard(d_i, d_j) = \frac{|G_{d_i} \cap G_{d_j}|}{|G_{d_i} \cup G_{d_j}|}$  (10)

201 where  $G_{d_i}$  and  $G_{d_j}$  denote the target gene set of  $d_i$  and  $d_j$ , respectively. Given a  
202 threshold  $\xi$ , we further estimate the percentage of drug pairs whose Jaccard indices  
203 exceed  $\xi$  as follows.

204  $Sim_U = \frac{| \{(d_i, d_j) | Jaccard(d_i, d_j) \geq \xi, (d_i, d_j) \in U \} |}{|U|}$  (11)

205 where  $U$  denotes the set of drug pairs and  $Sim_U$  measures the similarities within the  
206 set of drug pairs  $U$ . If  $\xi = \min_{(d_i, d_j) \in U} \frac{1}{|G_{d_i} \cup G_{d_j}|}$ , then  $Sim_U$  measures the  
207 percentage of drug pairs that target at least one common gene.

208 Two drugs also potentially interact through target genes that communicate via  
209 physical paths in human PPI networks. Given a gene pair  $(g_i, g_j)$ , we use breadth-first  
210 graph search algorithm to obtain all the physical paths between  $g_i$  and  $g_j$  from

211 human PPI networks, denotes as  $P_{(g_i, g_j)}$ . The length of the shortest path and longest  
 212 path within  $P_{(g_i, g_j)}$  is denoted as  $S_{(g_i, g_j)}$  and  $L_{(g_i, g_j)}$ , respectively. Then the average  
 213 number of paths  $Avg_{(d_i, d_j)}$ , the shortest path length  $S_{(d_i, d_j)}$  and the longest path  
 214 length  $L_{(d_i, d_j)}$  between drug  $d_i$  and  $d_j$  are defined as follows.

$$\begin{aligned}
 215 \quad Avg_{(d_i, d_j)} &= \frac{\sum_{(g_i, g_j), g_i \in G_{d_i} \wedge g_j \in G_{d_j}} |P_{(g_i, g_j)}|}{|\{(g_i, g_j) | g_i \in G_{d_i} \wedge g_j \in G_{d_j}\}|} \\
 216 \quad S_{(d_i, d_j)} &= \min_{(g_i, g_j), g_i \in G_{d_i} \wedge g_j \in G_{d_j}} S_{(g_i, g_j)} \\
 217 \quad L_{(d_i, d_j)} &= \max_{(g_i, g_j), g_i \in G_{d_i} \wedge g_j \in G_{d_j}} L_{(g_i, g_j)} \quad (12)
 \end{aligned}$$

218 These three metrics could measure the communication intensities between two drugs'  
 219 target genes. Especially,  $S_{(d_i, d_j)} = 0$  indicates that there exist common target genes  
 220 between drug  $d_i$  and  $d_j$ , and  $Avg_{(d_i, d_j)} = 0$  indicates that no paths exist between the  
 221 target genes of drug  $d_i$  and  $d_j$ .

## 222 Results

### 223 Performance of cross validation and independent test

224 The ROC curve of 5-fold cross validation is illustrated in Figure 1A. The performance  
 225 is fairly encouraging with ROC-AUC score equal to 0.9877. The other metrics for 5-  
 226 fold cross validation are provided in Table 1. The metrics of SP, SE and MCC show  
 227 that the proposed framework shows little bias between the positive and negative classes,  
 228 e.g. 0.9538 and 0.9394 SE on the positive and negative class, respectively. The overall  
 229 MMC is up to 0.8983. These results show that simple drug target profile is sufficient to  
 230 separate interacting drug pairs from non-interacting drug pairs with a high accuracy  
 231 (Acc=94.66%). Exploration of target genes that two drugs act on helps to clearly  
 232 unravel molecular mechanisms behind drug-drug interactions and reduce the

233 information chaos resulting from integration of heterogeneous data. Especially, the  
234 information of drug molecule structure is not indispensable and the less informative or  
235 irrelevant information would not be introduced into the proposed framework.

236 To estimate how well the proposed framework generalizes to unseen examples, we  
237 further conduct independent test on 13 external DDI datasets and one negative  
238 independent test data that are disjoint with the training data. The independent test data  
239 sizes are illustrated in Figure 1B varying from 8188 to 3. The recall rates are illustrated  
240 in Figure 1C. All the independent test data achieve no lower than 0.8 recall rate except  
241 “DDI Corpus 2013”. For instance, the proposed framework correctly recognizes  
242 94.97%, 89.92% and 97.30% of the experimental DDIs from KEGG [26], OSCAR [27]  
243 and VA NDF-RT [28], respectively (see Table 1). In addition, the proposed framework  
244 achieves 0.9373 recall rate on the negative independent test data, indicating a low risk  
245 of model bias. The independent test performance shows that the proposed framework  
246 trained via simple drug target profile could well generalize to unseen DDIs and  
247 meanwhile indicate that the encouraging cross validation performance does not result  
248 from overtraining.

## 249 **Molecular mechanism analyses of drug-drug interactions**

250 *Jaccard index between two drugs.* The overlap between two drugs’ target genes  
251 measured via Formula (11) could reveal some common molecular mechanisms that the  
252 drugs take effect. The comparison of Jaccard index between the positive training data  
253 and the negative training data is illustrated in Figure 1. The threshold of Jaccard index  
254 assumes  $\xi = \min_{(d_i, d_j) \in U} \frac{1}{|G_{d_i} \cup G_{d_j}|}$  and  $\xi = 0.5$  in Figure 2A and Figure 2B,  
255 respectively. The results show that interacting drugs tend to target more common genes.

256 *Average number of paths between two drugs.* The statistics of paths between two  
257 drugs' target genes in PPI networks measured via Formula (12) could reveal the  
258 interaction intensity and efficiency between drugs. To reduce the time of paths search,  
259 we only randomly choose 9692 interacting drug pairs and 9692 non-interacting drug  
260 pairs for molecular mechanism analyses. The average number of paths of top twenty  
261 drug pairs are illustrated in Figure 3A. We can see that interacting drug pairs have more  
262 paths between target genes than non-interacting drug pairs, indicating that heavy  
263 communication traffic between two drugs' target genes tends to cause drug-drug  
264 interactions. Additionally, there are fewer interacting drug pairs that no paths exist  
265 between their target genes than non-interacting drug pairs as shown in Figure 3B.

266 *Shortest path length between two drugs.* Computational results show that the length  
267 of shortest paths between two drugs' target genes ranges from 0 to 5. The statistics of  
268 shortest path length between two drugs is illustrated in Figure 3C. We can see that  
269 interacting drug pairs significantly outnumber non-interacting drug pairs in the case that  
270 the shortest path length is equal to 0 (i.e. common target genes). With the increase of  
271 shortest path length, non-interacting drug pairs gradually outnumber interacting drug  
272 pairs. These results show that common target genes or target genes communicating via  
273 shorter shortest physical paths tend to cause drug-drug interactions.

274 *Longest path length between two drugs.* Computational results show that the length  
275 of longest paths between two drugs' target genes ranges from 0 to 8. As shown in Figure  
276 3D, non-interacting drug pairs outnumber interacting drug pairs when the longest path  
277 ranges from 3 to 5, but conversely interacting drug pairs significantly outnumber non-  
278 interacting drug pairs when the longest path length equals to 6. These results to some  
279 extent show that interacting drugs potentially act on each other via a long range of effect  
280 propagation. The metrics  $Avg_{(d_i,d_j)}$ ,  $S_{(d_i,d_j)}$  and  $L_{(d_i,d_j)}$  as defined in Formula (12)

281 could measure the tendency of drug-drug interaction from the aspects of interaction  
282 intensity, interaction efficacy and action range. The shortest path length equal to 0 and  
283 the longest path length equal to 6 are the best indicators to distinguish between  
284 interacting drug pairs and non-interacting drug pairs.

285 *Common target pathways between two drugs.* We map the target genes onto NetPath  
286 [33] and Reactome [34] signaling pathways to gain knowledge about the number of  
287 common target pathways between two drugs. Computational results show that  
288 interacting drug pairs tend to target more common signaling pathways than non-  
289 interacting drug pairs (see Figure 4A for NetPath pathways and Figure 4B for Reactome  
290 pathways). Targeting genes located at the same signaling pathways helps to induce  
291 synergistic or antagonistic effects of one drug to other drugs.

292 *Common cellular processes between two drugs.* Similar to the statistics of common  
293 target pathways, the target genes of two interacting drugs are more likely to be involved  
294 in common cellular processes than those of two non-interacting drugs (see Figure 4C).  
295 This phenomenon is easily interpreted. Two drugs that target genes participating  
296 common cellular processes more likely alter each other's therapeutic effects.

## 297 **Comparison with existing methods**

298 The fundamental difference between this proposed framework and existing methods is  
299 that this framework assumes the reactions between target genes of two drugs to be the  
300 major driving factor causing drug-drug interactions, while existing methods majorly  
301 focus on inferring drug-drug interactions from drug structural similarities. From  
302 computational point of view, this framework only uses simple representation of drug  
303 target profiles to depict drug pairs so that the model is easy to interpret and the model  
304 complexity is under control. However, existing methods integrate multiple  
305 heterogeneous data so that we do not know which data is informative and how much

306 the data contributes to the performance. More importantly, the model complexity is  
307 increased and the model is hard to interpret.

308 From quantitative point of view, we compare this proposed framework with seven  
309 existing methods in terms of performance as provided in Table 2. Existing methods  
310 generally achieve ROC-AUC scores except Cheng et al. [15] (ROC-AUC=0.67).  
311 Nevertheless, these methods generally show imbalanced performance on the positive  
312 and negative class. For instance, Vilar et al. [6] use drug structure profile to infer drug-  
313 drug interactions but shows extreme bias towards the negative class (e.g. SE 0.68 and  
314 0.96 for the positive and negative class, respectively). Zhang et al. [16] integrate  
315 heterogeneous data of drug substructures, drug targets, drug enzymes, drug transporters,  
316 drug pathways, drug indications and drug side-effects data, achieving encouraging  
317 performance (ROC-AUC score=0.957, PR=0.785, SE=0.670) but only recognizing 7  
318 out of 20 predicted DDIs (equivalent to 35% recall rate of independent test) . Similarly,  
319 Gottlieb et al. [20] achieve fairly good performance of cross validation but achieve only  
320 53% recall rate of independent test.

321 Deep learning has been used to predict effects or types of drug-drug interaction [18,  
322 19]. Karim et al. [22] use deep learning to learn the representation of DDI networks  
323 structures and predict novel DDIs. This method also achieves satisfactory performance  
324 (ROC-AUC score=0.97, MCC=0.79, F1 score=0.91) but is hard to interpret from the  
325 molecular action mechanisms between two drugs. Comparing Table 1 and Table 2, we  
326 can see that this proposed framework also outperforms existing methods in terms of  
327 performance.

## 328 **Predictions and clinical implications**

329 We randomly sample 99,986 drug pairs disjoint with the training data and independent  
330 test data as prediction set. The proposed framework predicts 43,719 drug pairs to

331 interact (see Supplementary file S1). These predictions potentially contain a certain  
332 level of false interactions. As each prediction is provided with a probability of  
333 confidence level, a higher threshold (e.g. 0.7) could be used to filter out those predicted  
334 weak interactions. For the predicted DDIs, we further analyse the common cellular  
335 processes (see Supplementary file S2) and common signaling pathways (see  
336 Supplementary file S3) that the drug target genes are involved in. These analyses help  
337 us to understand the underlying mechanisms behind drug-drug interactions. As a case  
338 study, we analyse the predicted interaction between drug Nabiximols and Glucosamine  
339 from the aspects of common cellular processes and signaling pathways.

#### 340 **Case study on predicted drug-drug interaction (Nabiximols, Glucosamine)**

341 Nabiximols ( $C_{42}H_{60}O_4$ ) extracted from *Cannabis sativa L.* is used for the treatment of  
342 neuropathic pain from Multiple Sclerosis and for intractable cancer pain, with the  
343 pharmacological effects of analgesic, muscle relaxant, anxiolytic, neuroprotective and  
344 anti-psychotic activity (<https://www.drugbank.ca/drugs/DB14011>). Glucosamine  
345 ( $C_6H_{13}NO_5$ ), as a precursor for glycosaminoglycans that are a major component of joint  
346 cartilage, is commonly used to rebuild cartilage and treat osteoarthritis  
347 (<https://www.drugbank.ca/drugs/DB01296>). According to DrugBank [24], Nabiximols  
348 targets 57 human genes and Glucosamine targets 6 human genes. We further analyse  
349 the common cellular processes and signaling pathways that the two drugs' target genes  
350 get involved in.

351 *Common cellular processes between Nabiximols and Glucosamine.* Computational  
352 results show that there are 68 common cellular processes that the target genes of  
353 Nabiximols and Glucosamine are involved in. For clarity, we only illustrate 21 cellular  
354 processes and their associated target genes in Figure 5 and the rest cellular processes  
355 are provided in Supplementary file S2. Two drugs mediate common cellular processes

356 via common target genes or two associated target genes involved in the same cellular  
357 processes. As shown in Figure 5, Nabiximols and Glucosamine mediate the common  
358 cellular processes of exogenous drug catabolic process (GO:0042738) and drug  
359 metabolic process (GO:0017144) via the common gene *CYP2C19*. In addition,  
360 Nabiximols and Glucosamine mediate the common cellular processes of negative  
361 regulation of smooth muscle cell proliferation (GO:0048662) via Nabiximols target  
362 gene *PPARG* and Glucosamine target gene *IFNG*. Take another example, Nabiximols  
363 and Glucosamine mediate the common cellular processes of regulation of reactive  
364 oxygen species (ROS) metabolic process (GO:2000377) via Nabiximols target gene  
365 *CYP1B1* and Glucosamine target gene *TNF*. Association via different target genes is  
366 one major way of two drugs mediating common cellular processes, which potentially  
367 leads to synergistic or antagonistic drug-drug interactions. The proposed framework  
368 predicts many DDIs that target no common genes but mediate common cellular  
369 processes via different target genes (see Supplementary file S2). For instance, drug  
370 Nabiximols (DB14011) and Gallium nitrate (DB05260) are not found to target common  
371 human genes at present, but they are predicted to target the common cellular processes  
372 of neutrophil chemotaxis (GO:0030593), positive regulation of NF-kappaB  
373 transcription factor activity (GO:0051092), etc.

374 *Common signaling pathways between Nabiximols and Glucosamine.* The common  
375 Reactome signaling pathways that Nabiximols and Glucosamine mediate via common  
376 or associated target genes are illustrated in Figure 6. Among the target genes, the  
377 common target gene *CYP2C19* is associated with four Reactome signaling pathways,  
378 i.e. Synthesis of epoxy (EET) and dihydroxyeicosatrienoic acids (DHET) (R-HSA-  
379 2142670), Xenobiotics (R-HSA-211981), CYP2E1 reactions (R-HSA-211999) and  
380 Synthesis of (16-20)-hydroxyeicosatetraenoic acids (HETE) (R-HSA-2142816).

381 Nabiximols and Glucosamine mediate the common signaling pathway of Neutrophil  
382 degranulation (R-HSA-6798695) via Nabiximols target gene *ALOX5* and Glucosamine  
383 target gene *MMP9*. Mediation of common signaling pathways also potentially leads to  
384 drug-drug interactions. Two drugs that do not target common human genes also  
385 potentially mediate the same signaling pathways (see Supplementary file S3). For  
386 instance, drug Nabiximols (DB14011) and SF1126 (DB05210) have not been reported  
387 to target common human genes, but they are predicted to mediate some common  
388 signaling pathways, e.g. Regulation of PTEN gene transcription (R-HSA-8943724),  
389 Interleukin-4 and Interleukin-13 signaling (R-HSA-6785807), G alpha (q) signaling  
390 events (R-HSA-416476).

## 391 **Discussion**

392 Drug-drug interactions are usually detected and reported only after co-prescribed drugs  
393 have clinically done damages to patient health and life. We need resort to computational  
394 methods to predict whether two drugs interact before clinical co-prescription to reduce  
395 potential risk of side effects. Existing methods generally explore multiple  
396 heterogeneous data sources via data integration to increase model performance. The  
397 major drawback of these methods lies in the high complexity of model and data. In  
398 these methods, we do not know which information contributes to the model  
399 performance and interprets the molecular mechanisms behind drug-drug interactions.  
400 Among the heterogeneous data, existing methods heavily depend on drug structures  
401 with the assumption that structurally similar drugs take into similar therapeutic effects.  
402 This assumption surely captures a fraction of drug-drug interactions but shows bias,  
403 because it ignores a large fraction of interactions between structurally dissimilar drugs.

404 Furthermore, data integration would fail when required data are not available, e.g. drug  
405 structures, drug side-effects, clinical records.

406 In this study, we propose a simple representation of drug target profiles to depict drug  
407 pairs, based on which to train an  $l_2$ -regularized logistic regression model for DDI  
408 prediction. As compared to the existing methods, this proposed method could directly  
409 interpret the molecular mechanisms behind drug-drug interactions via the association  
410 and action between two drugs' target genes in the context of protein-protein interactions  
411 (PPI) networks and signaling pathways. We use the known drug-drug interactions from  
412 DrugBank as the positive training data and randomly sample the same size drug pairs  
413 as the negative training data to train an  $l_2$ -regularized logistic regression model.  
414 Computational results show that the proposed framework achieves fairly encouraging  
415 performance of cross validation and outperforms the existing methods. Furthermore,  
416 the proposed framework demonstrates fairly good performance of independent test on  
417 thirteen external DDI datasets. The encouraging performance on the randomly sampled  
418 negative independent test data demonstrates no bias towards the positive class.  
419 However, the proposed framework demonstrates a little large fraction of false  
420 interactions on the prediction set. This result is largely due to the quality of randomly  
421 sampled negative training data. Lack of experimentally verified non-interacting drug  
422 pairs is a major concern to be addressed in computational modelling. Here we adopt a  
423 higher threshold of probability to filter out the weak predictions.

424 In addition, we develop three statistical metrics to measure the communication  
425 intensity, interaction efficacy and action range between two drugs. These metrics help  
426 us to understand the statistical characteristic of the paths between two drugs' target  
427 genes in human PPI networks. In addition, we use common cellular processes and  
428 signaling pathways between two drugs to understand the mechanisms behind drug-drug

429 interactions. The unravelled mechanisms provide biological insights into potential  
430 pharmacological risks of known DDIs and drug target genes.

## 431 **Conclusions**

432 Integration of heterogeneous data is not indispensable for drug-drug interaction  
433 prediction. Simple representation of drug target profile alone achieves encouraging  
434 model performance that outperforms existing methods with intuitive elucidation of  
435 molecular mechanisms of drug-drug interactions.

## 436 **Additional Information**

437 **File S1** Text file contains the predicted drug-drug interactions.  
438 (TXT)

439 **File S2** Text file contains the genes that are predicted to be targeted by novel drugs  
440 for Matador.  
441 (TXT)

442 **File S3** Text file contains the genes that are predicted to be targeted by novel drugs  
443 for DrugBank.  
444 (TXT)

## 445 **Declarations**

## 446 **Statements of ethical approval**

447 Not applicable

## 448 **Acknowledgement**

449 This work is partly supported by the funding from the NIH grants 2U54MD007595,  
450 5P20GM103424-17 and U19AG055373. The contents are solely the responsibility of  
451 the authors and do not necessarily represent the official views of the NIH.

## 452 **Author Contributions**

453 MS conducted the study and wrote the paper.

## 454 **Competing interests**

455 The author declares that there no competing interests.

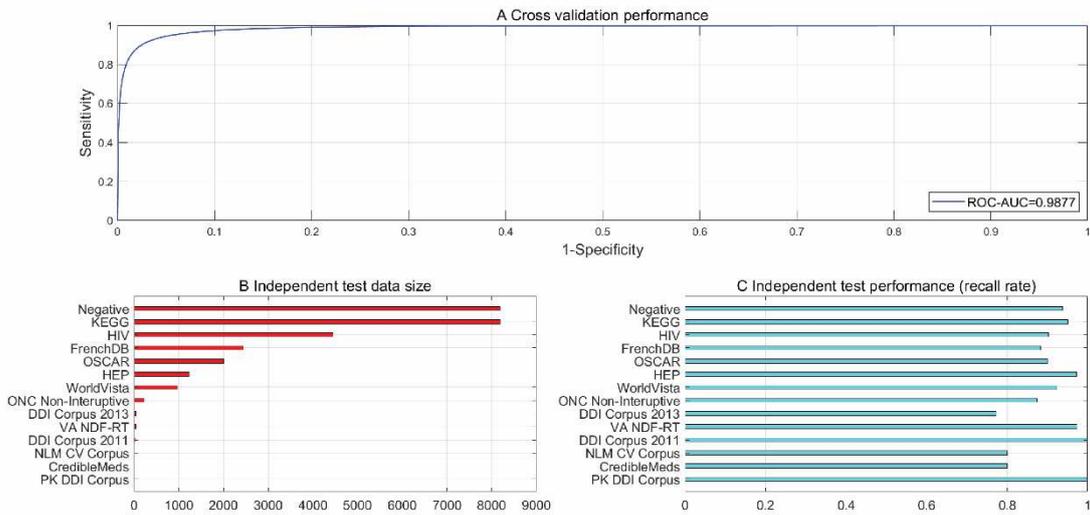
## 456 **References**

- 457 1. Wienkers LC, Heath TG. Predicting in vivo drug interactions from in vitro drug  
458 discovery data. *Nat Rev Drug Discovery* 4:825–33 (2005).
- 459 2. Edwards IR, Aronson JK. Adverse drug reactions: definitions, diagnosis, and  
460 management. *Lancet* 356:1255–9 (2000).
- 461 3. Leape LL, Bates DW, Cullen DJ, et al. Systems analysis of adverse drug events.  
462 ADE Prevention Study Group. *JAMA* 274:35–43 (1995).
- 463 4. Duke JD, Han X, Wang Z, Subhadarshini A, Karnik SD, et al. Literature based  
464 drug interaction prediction with clinical assessment using electronic medical  
465 records: novel myopathy associated drug interactions. *PLoS Comput Biol*  
466 8:e1002614 (2012).
- 467 5. Vilar S, Hripcsak G. The role of drug profiles as similarity metrics: applications  
468 to repurposing, adverse effects detection and drug-drug interactions. *Brief*  
469 *Bioinform* 18:670-681 (2017).
- 470 6. Vilar S, Harpaz R, Uriarte E, Santana L, Rabadan R, et al. Drug-drug interaction  
471 through molecular structure similarity analysis. *J Am Med Inform Assoc*  
472 19:1066-74 (2012).
- 473 7. Ferdousi R, Safdari R, Omid Y. Computational prediction of drug-drug  
474 interactions based on drugs functional similarities. *J Biomed Inform* 70:54-64  
475 (2017).
- 476 8. Vilar S, Uriarte E, Santana L, Lorberbaum T, Hripcsak G et al. Similarity-based  
477 modeling in large-scale prediction of drug-drug interactions. *Nat Protoc* 9:  
478 2147–2163 (2014).
- 479 9. Zhang W, Chen Y, Li D, Yue X. Manifold regularized matrix factorization for  
480 drug-drug interaction prediction. *J Biomed Inform* 88:90-97 (2018).
- 481 10. Shtar G, Rokach L, Shapira B. Detecting drug-drug interactions using artificial  
482 neural networks and classic graph similarity measures. *PLoS One* 14:e0219796  
483 (2019).
- 484 11. Zhang P, Wang F, Hu J, Sorrentino R. Label Propagation Prediction of Drug-  
485 Drug Interactions Based on Clinical Side Effects. *Sci Rep* 5:12339 (2015).
- 486 12. Park K, Kim D, Ha S, Lee D. Predicting Pharmacodynamic Drug-Drug  
487 Interactions through Signaling Propagation Interference on Protein-Protein  
488 Interaction Networks. *PLoS One* 10:e0140816 (2015).
- 489 13. Huang J, Niu C, Green CD, Yang L, Mei H, Han JD. Systematic prediction of  
490 pharmacodynamic drug-drug interactions through protein-protein-interaction  
491 network. *PLoS Comput Biol* 9:e1002998 (2013).
- 492 14. Dhami DS, Kunapuli G, Das M, Page D, Natarajan S. Drug-Drug Interaction  
493 Discovery: Kernel Learning from Heterogeneous Similarities. *Smart Health*  
494 (Amst) 9-10:88-100 (2018).
- 495 15. Cheng F, Zhao Z. Machine learning-based prediction of drug-drug interactions  
496 by integrating drug phenotypic, therapeutic, chemical, and genomic properties.  
497 *J Am Med Inform Assoc* 21:e278-86 (2014).

- 498 16. Zhang W, Chen Y, Liu F, Luo F, Tian G, Li X. Predicting potential drug-drug  
499 interactions by integrating chemical, biological, phenotypic and network data.  
500 BMC Bioinformatics 18:18 (2017).
- 501 17. Song D, Chen Y, Min Q, Sun Q, Ye K, Zhou C, et al. Similarity-based machine  
502 learning support vector machine predictor of drug-drug interactions with  
503 improved accuracies. J Clin Pharm Ther 44:268-275 (2019).
- 504 18. Ryu JY, Kim HU, Lee SY. Deep learning improves prediction of drug-drug and  
505 drug-food interactions. Proc Natl Acad Sci U S A 115:E4304-E4311 (2018).
- 506 19. Lee G, Park C, Ahn J. Novel deep learning model for more accurate prediction  
507 of drug-drug interaction effects. BMC Bioinformatics 20:415 (2019).
- 508 20. Gottlieb A, Stein GY, Oron Y, Ruppin E, Sharan R. INDI: a computational  
509 framework for inferring drug interactions and their associated  
510 recommendations. Mol Syst Biol 8:592 (2012).
- 511 21. Qian S, Liang S, Yu H. Leveraging genetic interactions for adverse drug-drug  
512 interaction prediction. PLoS Comput Biol 15:e1007068 (2019).
- 513 22. Karim MR, Cochez M, Jares JB, Uddin M, Beyan O, Decker S. Drug-Drug  
514 Interaction Prediction Based on Knowledge Graph Embeddings and  
515 Convolutional-LSTM Network. arXiv:1908.01288 (2019).
- 516 23. Jia J, Zhu F, Ma X, Cao Z, Cao ZW, et al. Mechanisms of Drug Combinations:  
517 Interaction and Network Perspectives. Nat Rev Drug Discov 8:111-28 (2009).
- 518 24. Wishart DS et al. DrugBank 5.0: a major update to the DrugBank database for  
519 2018. Nucleic Acids Res 46:D1074-D1082 (2018).
- 520 25. Ayvaz S et al. Toward a Complete Dataset of Drug-Drug Interaction  
521 Information From Publicly Available Sources. J Biomed Inform 55: 206-17  
522 (2015).
- 523 26. Kanehisa M, Goto S, Sato Y, Kawashima M, Furumichi M et al. Data,  
524 information, knowledge and principle: back to metabolism in KEGG. Nucleic  
525 Acids Res 42(Database issue): D199–D205 (2014).
- 526 27. Crowther NR, Holbrook AM, Kenwright R, Kenwright M. Drug interactions  
527 among commonly used medications. Chart simplifies data from critical  
528 literature review. Can Fam Physician 43:1972-6, 1979-81 (1997).
- 529 28. Olvey EL, Clauschee S, Malone DC. Comparison of critical drug–drug  
530 interaction listings: the department of Veterans Affairs medical system and  
531 standard reference compendia. Clin Pharmacol Ther 87:48–51 (2010).
- 532 29. Keshava Prasad TS, Goel R, Kandasamy K, Keerthikumar S, Kumar S et al.  
533 Human Protein Reference Database--2009 update. Nucleic Acids Res  
534 37(Database issue), D767-72 (2009).
- 535 30. Chatr-Aryamontri A, Breitkreutz BJ, Oughtred R, Boucher L, Heinicke S et al.  
536 The BioGRID interaction database: 2015 update. Nucleic Acids Res 43  
537 (Database issue), D470-8 (2015).
- 538 31. Orchard S, Ammari M, Aranda B, Breuza L, Briganti L. The MIntAct project–  
539 IntAct as a common curation platform for 11 molecular interaction databases.  
540 Nucleic Acids Res. (Database issue) 42, D358–63 (2014).

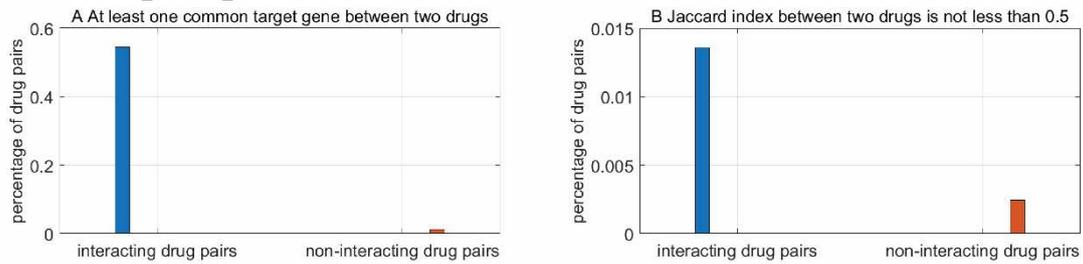
541 32. López Y, Nakai K, Patil A. HitPredict version 4: comprehensive reliability  
542 scoring of physical protein-protein interactions from more than 100 species.  
543 Database (Oxford) pii: bav117 (2015).  
544 33. Kandasamy K, et al. NetPath: a public resource of curated signal transduction  
545 pathways. Genome Biol 11: R3 (2010).  
546 34. Fabregat A, et al. The Reactome Pathway Knowledgebase. Nucleic Acids  
547 Res46(Database issue): D649–D655 (2018).  
548 35. Fan, R, Chang, K, Hsieh, C, Wang, X & Lin, C. LIBLINEAR: A Library for  
549 Large Linear Classification. Mach Learn Res 9, 1871-1874 (2008).  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582

583 **Figure 1 Performance of cross validation and independent test. A.** ROC  
 584 curve and AUC score for 5-fold cross validation. **B.** Statistics of independent  
 585 test data size. **C.** Recall rates on the independent test data.



586  
 587  
 588  
 589  
 590  
 591  
 592  
 593  
 594  
 595  
 596  
 597  
 598  
 599  
 600  
 601  
 602  
 603  
 604  
 605  
 606  
 607  
 608  
 609  
 610  
 611  
 612  
 613  
 614  
 615  
 616  
 617

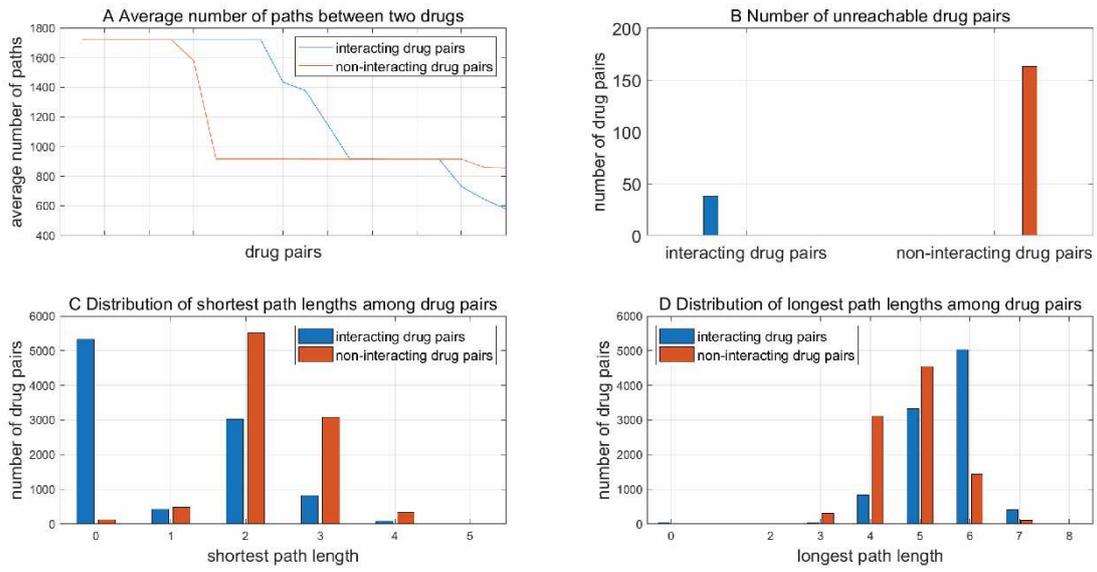
618 **Figure 2 Statistics of common target genes between interacting and non-**  
619 **interacting drugs.**



620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661

662  
663

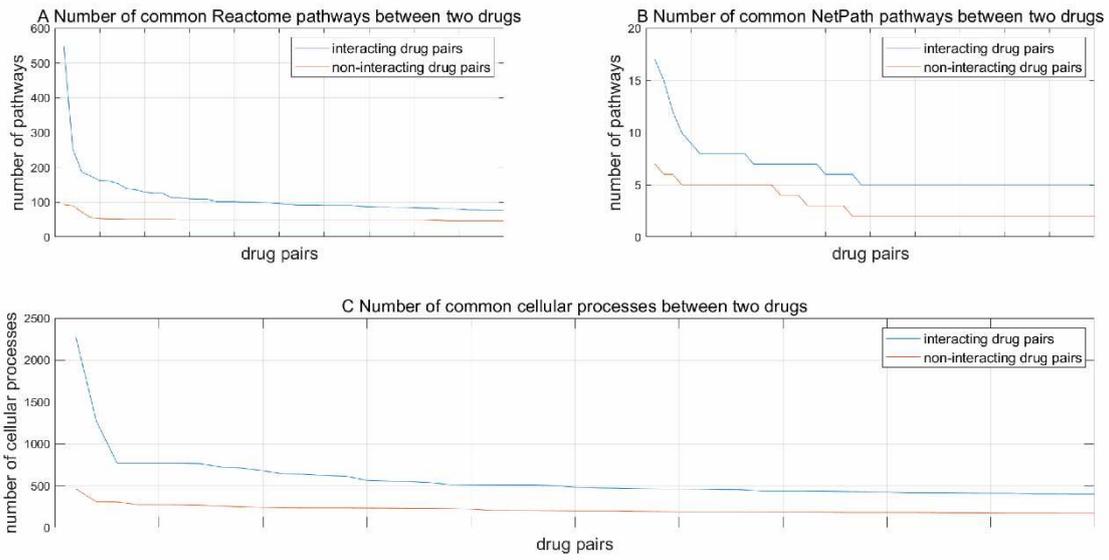
**Figure 3 The statistics of average number of paths, shortest path lengths and longest path lengths between two drugs.**



664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696

697  
698

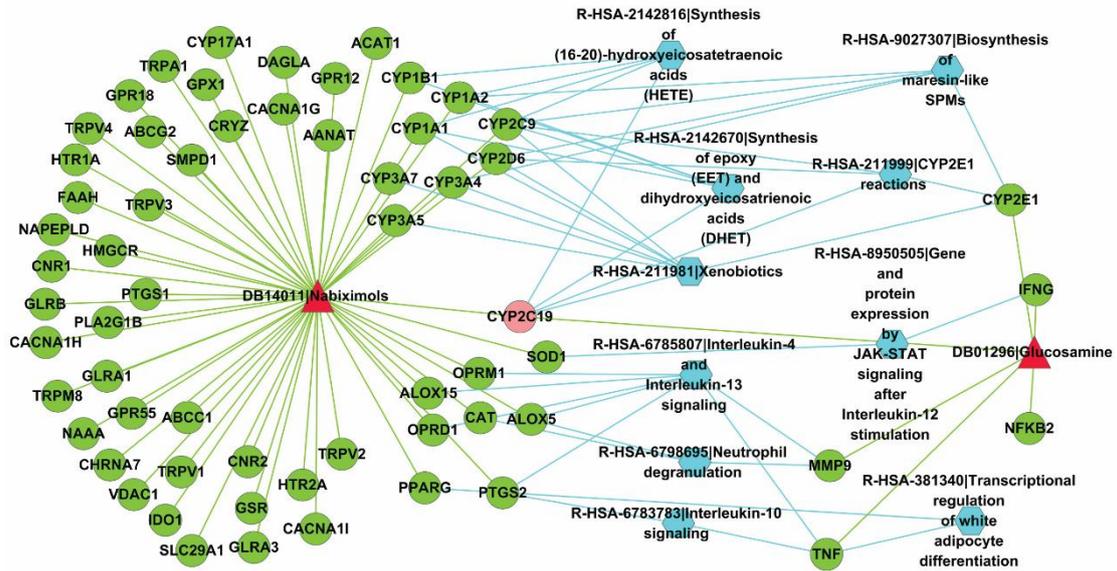
**Figure 4 Statistics of common signaling pathways that two drugs target and common cellular processes that two drugs are involved in.**



699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731



756 **Figure 6 Common target Reactome signaling pathways between**  
 757 **DB14011|Nabiximols and DB01296|Glucosamine predicted to interact.**  
 758 Red triangle nodes denote drugs; green circle nodes denote drug target genes;  
 759 light red circle nodes denote common target genes; and blue hexagon nodes  
 760 denote Reactome signaling pathways.



761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776 **Table 1 Performance estimation of 5-fold cross validation and**  
 777 **independent test.**

Cross validation							Independent test (recall rate)			
PR	SE	MCC	Acc	MCC*	AUC	F1 score	KEGG	OSCAR	VA NDF-RT	Negative
0.9402(+)	0.9538(+)	0.8985(+)	94.66%	0.8983	0.9877	0.9470	0.9497	0.8992	0.9730	0.9373
0.9531(-)	0.9394(-)	0.8983(-)								

Note: + denotes positive class, - denotes negative class and MCC\* denotes overall MCC.

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

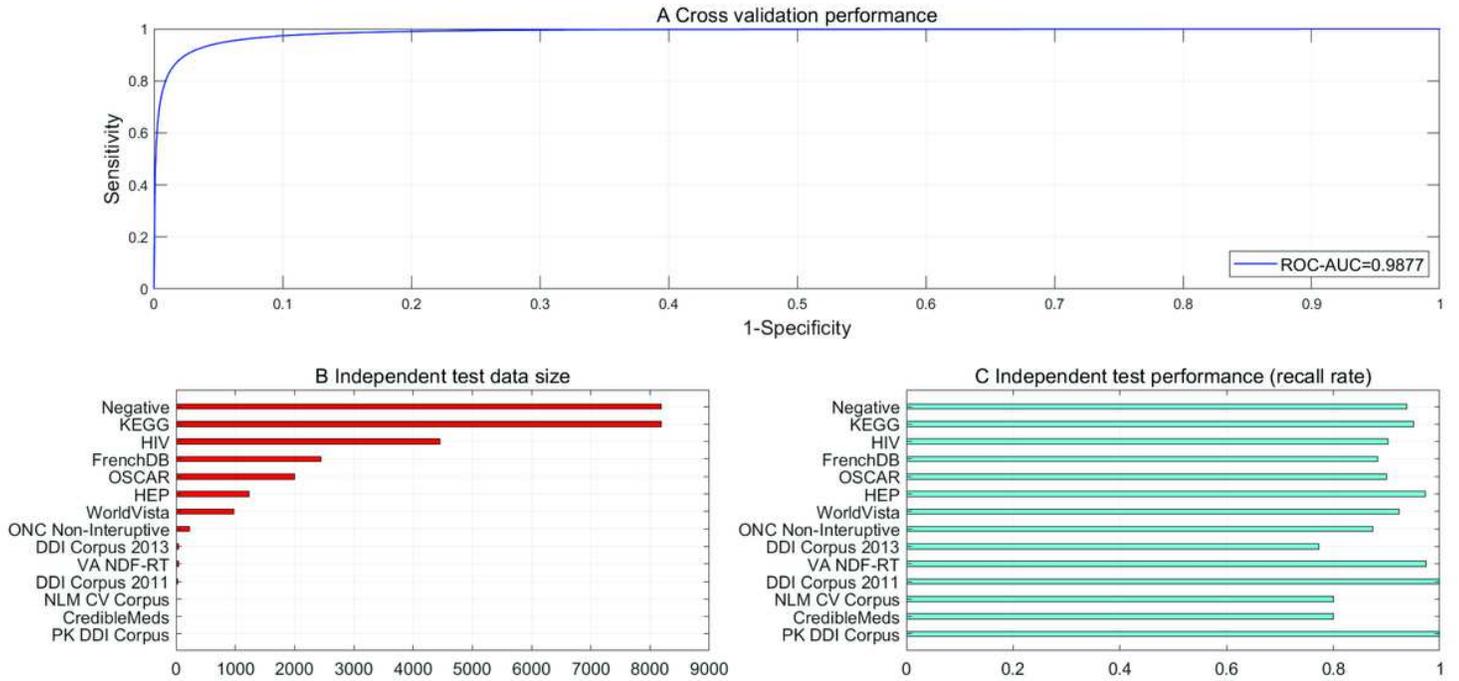
797

**Table 2 Performance comparisons with existing methods.**

	Cross validation		MCC	F1 score	ROC-AUC	Independent test
	PR	SE				
Vilar et al. [6]	0.26 (+) 11.81(-)	0.68 (+) 0.96(-)	-	-	0.92	31%
Ferdousi et al. [7]	-	0.72(+)	-	-	-	-
Cheng et al. [15]	-	-	-	-	0.67	-
Zhang et al. [16]	0.785	0.670	-	0.723	0.957	35%
Song et al. [17]	0.68 (+)	-	-	-	0.9738	24%
Gottlieb et al. [20]	0.88	0.93	-	-	0.96	53%
Karim et al. [22]	-	-	0.79	0.91	0.97	-

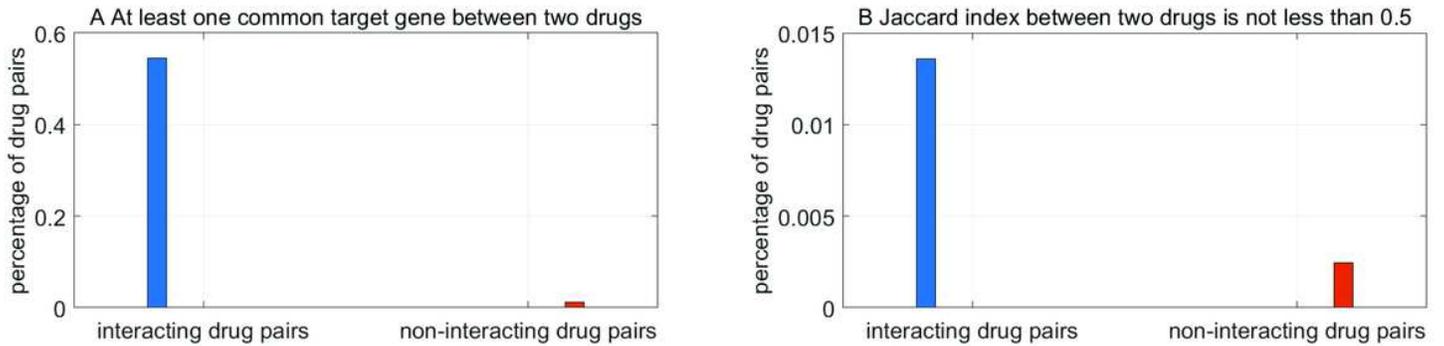
Note: + denotes positive class and - denotes negative class.

# Figures



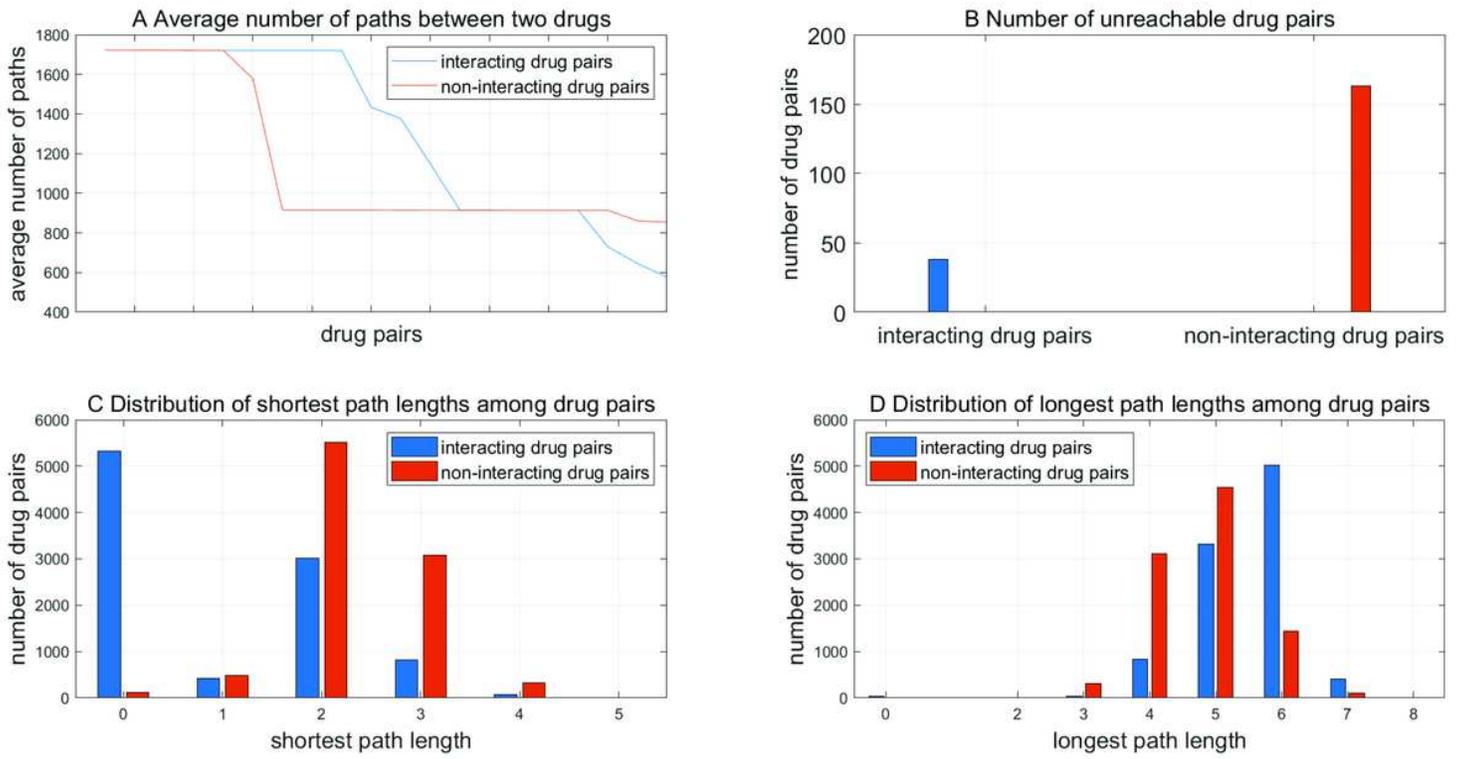
**Figure 1**

Performance of cross validation and independent test. A. ROC curve and AUC score for 5-fold cross validation. B. Statistics of independent test data size. C. Recall rates on the independent test data.



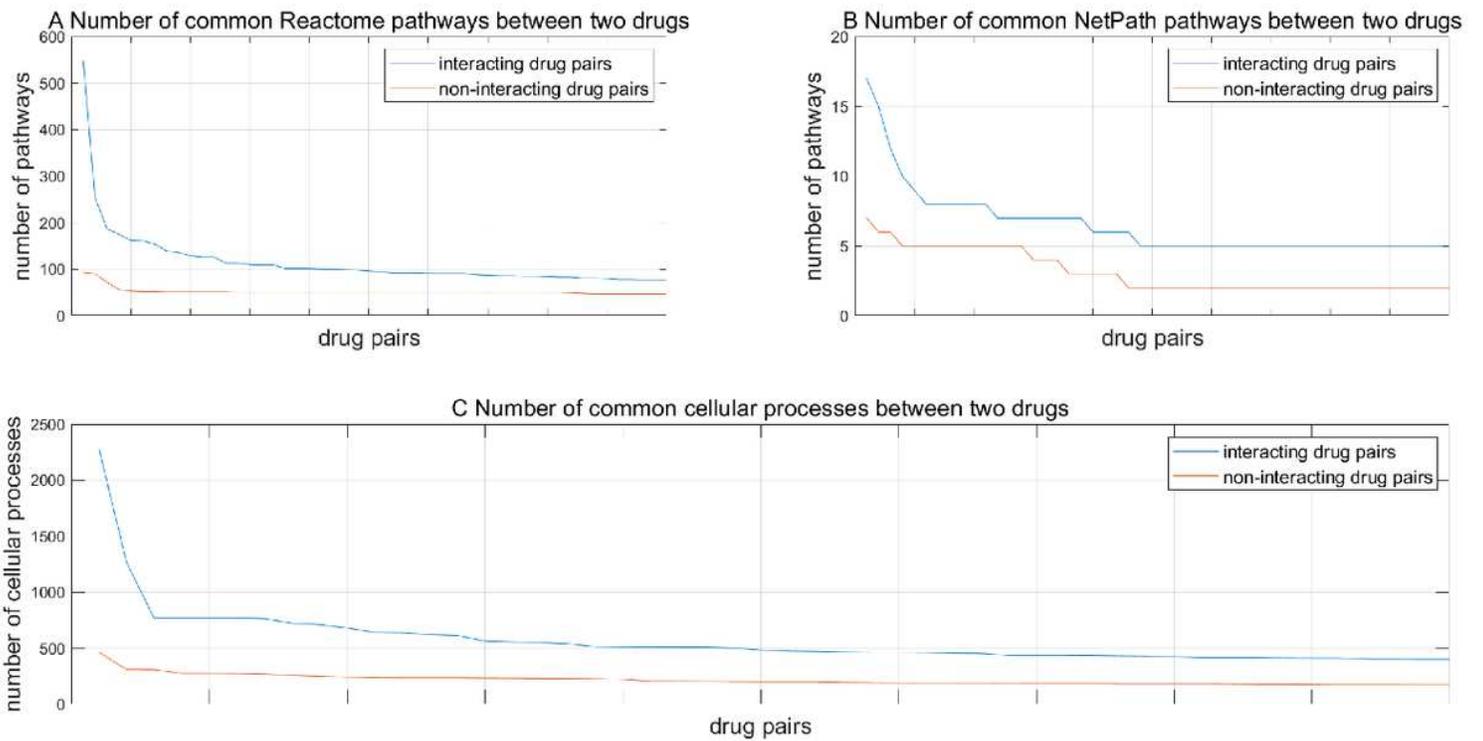
**Figure 2**

Statistics of common target genes between interacting and non-interacting drugs.



**Figure 3**

The statistics of average number of paths, shortest path lengths and longest path lengths between two drugs.



**Figure 4**



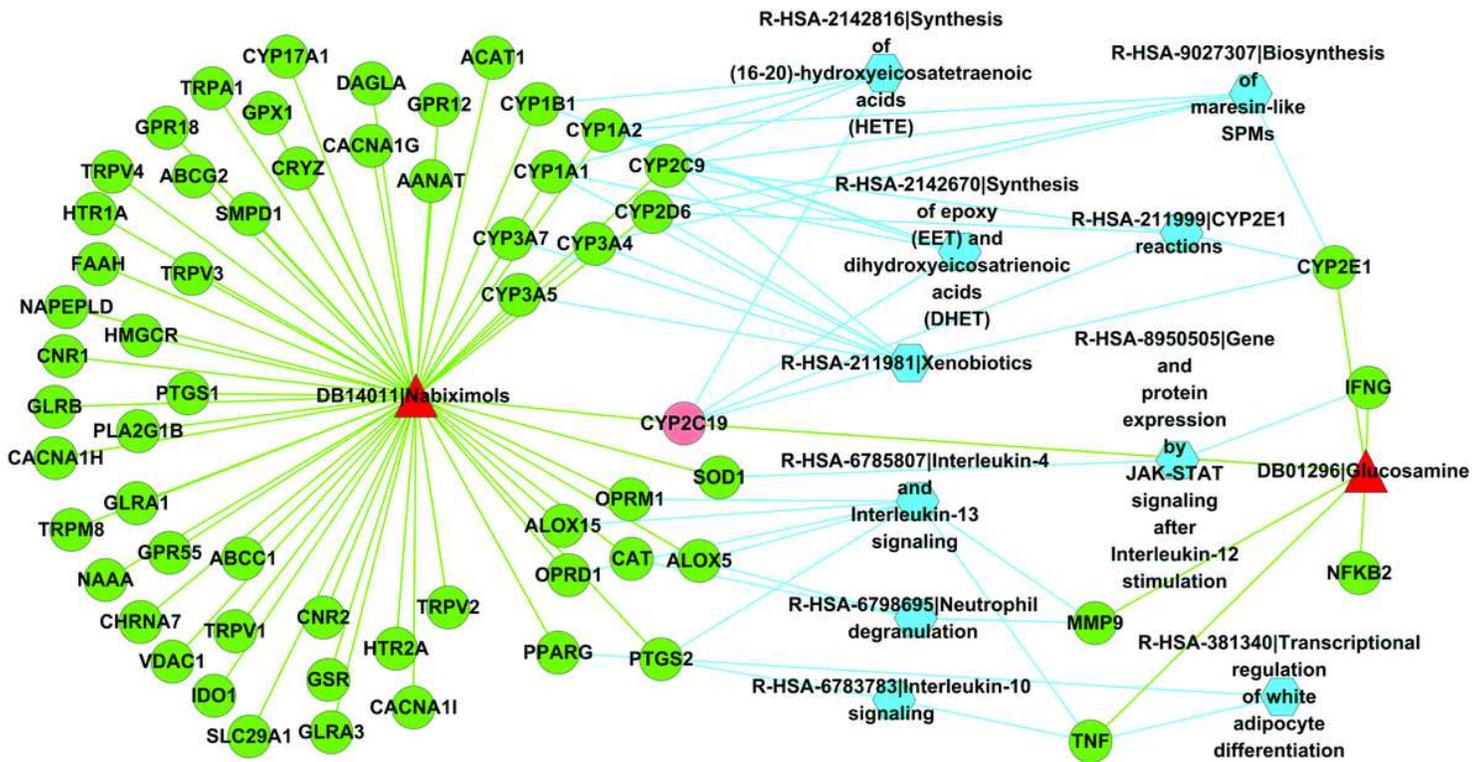


Figure 6

Common target Reactome signaling pathways between DB14011|Nabiximols and DB01296|Glucosamine predicted to interact. Red triangle nodes denote drugs; green circle nodes denote drug target genes; light red circle nodes denote common target genes; and blue hexagon nodes denote Reactome signaling pathways.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [S1.txt](#)
- [S2.txt](#)
- [S3.txt](#)